

**IN THE UNITED STATES PATENT AND TRADEMARK OFFICE**

In re U.S. Patent Application of )  
KIMURA et al. )  
Application Number: To be assigned )  
Filed: Concurrently herewith )  
For: METHOD FOR INDICATING RELATIONSHIP )  
BETWEEN CDNA SEQUENCE AND GENOME )  
RECORDING MEDIUM, SEQUENCER APPARATUS, )  
AND METHOD FOR DESIGNING A PRIMER )

Honorable Assistant Commissioner  
for Patents  
Washington, D.C. 20231

**REQUEST FOR PRIORITY  
UNDER 35 U.S.C. § 119  
AND THE INTERNATIONAL CONVENTION**

Sir:

In the matter of the above-captioned application for a United States patent, notice is hereby given that the Applicant claims the priority date of September 25, 2000, the filing date of the corresponding Japanese patent application 2000-289728.

The certified copy of corresponding Japanese patent application 2000-289728 is submitted herewith. Acknowledgment of receipt of the certified copy is respectfully requested in due course.

Respectfully submitted,

\_\_\_\_\_  
Stanley P. Fisher  
Registration Number 24,344

**REED SMITH HAZEL & THOMAS LLP**  
3110 Fairview Park Drive  
Suite 1400  
Falls Church, Virginia 22042  
(703) 641-4200  
August 21, 2001

j1046 U.S. PTO  
09/933168  
08/21/01

日 本 国 特 許 庁  
JAPAN PATENT OFFICE

J1046 U.S. PTO  
09/933168  
08/21/01

別紙添付の書類に記載されている事項は下記の出願書類に記載されている事項と同一であることを証明する。

This is to certify that the annexed is a true copy of the following application as filed with this Office

出 願 年 月 日

Date of Application:

2000年 9月25日

出 願 番 号

Application Number:

特願2000-289728

出 願 人

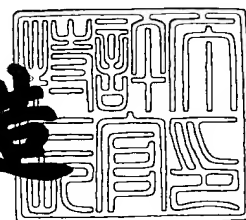
Applicant(s):

株式会社日立製作所

2001年 6月 1日

特 許 庁 長 官  
Commissioner,  
Japan Patent Office

及 川 耕 造



出証番号 出証特2001-3051656

【書類名】 特許願

【整理番号】 H001612

【提出日】 平成12年 9月25日

【あて先】 特許庁長官 殿

【国際特許分類】 G01N 33/50

【発明者】

    【住所又は居所】 東京都国分寺市東恋ヶ窪一丁目 2 8 0 番地 株式会社  
                                 日立製作所 中央研究所内

    【氏名】 木村 宏一

【発明者】

    【住所又は居所】 東京都国分寺市東恋ヶ窪一丁目 2 8 0 番地 株式会社  
                                 日立製作所 中央研究所内

    【氏名】 西川 哲夫

【特許出願人】

    【識別番号】 000005108

    【氏名又は名称】 株式会社 日立製作所

【代理人】

    【識別番号】 100091096

    【弁理士】

    【氏名又は名称】 平木 祐輔

【手数料の表示】

    【予納台帳番号】 015244

    【納付金額】 21,000円

【その他】 国等の委託研究の成果に係る特許出願（平成12年度新  
                                 エネルギー・産業技術総合開発機構（再）委託研究、産  
                                 業活力再生特別措置法第30条の適用を受けるもの）

【提出物件の目録】

    【物件名】 明細書 1

    【物件名】 図面 1

特 2 0 0 0 - 2 8 9 7 2 8

【物件名】	要約書	1
【プルーフの要否】	要	

【書類名】 明細書

【発明の名称】 cDNA配列とゲノム配列との対応表示方法、記録媒体、シーケンサ装置及びプライマ設計方法

【特許請求の範囲】

【請求項1】 グラフの1の軸にゲノム配列上の塩基位置を、他の軸にcDNA配列上の塩基位置をとり、前記ゲノム配列のうち所定塩基長以上を有する部分配列について、前記cDNA配列との間で所定の割合以上の類似性を有する部分をグラフ上に線分で表示することを特徴とするcDNA配列とゲノム配列との対応表示方法。

【請求項2】 複数のcDNAを縦軸にとり、前記cDNAとの対応関係をcDNAごとに異なる色で表示することを特徴とする請求項1記載のcDNA配列とゲノム配列との対応表示方法。

【請求項3】 ゲノム配列とcDNA配列とを入力するステップと、前記ゲノム配列中の所定の塩基長以上を有する部分配列について前記cDNA配列との間で所定の割合以上の類似性を有する部分を検索するステップと、前記ゲノム配列とcDNA配列をそれぞれグラフの縦軸と横軸又は横軸と縦軸にとって前記検索するステップで検索された部分を線分で該グラフ上に表示するステップと、を備えるcDNA配列とゲノム配列との対応表示方法をコンピュータに実行させるためのプログラムを記録したことを特徴とするコンピュータ読み取り可能な記録媒体。

【請求項4】 さらに前記所定の塩基長及び前記類似性の所定の割合を入力するステップを備えるcDNA配列とゲノム配列との対応表示方法をコンピュータに実行させるためのプログラムを記録したことを特徴とする請求項3記載の記録媒体。

【請求項5】 ネットワーク接続された又は内蔵されたゲノムデータベースにアクセスし、ゲノム配列を入力し、シーケンシングによって得られたcDNA配列を入力する入力手段と、前記ゲノム配列中の所定の塩基長以上を有する部分配列について前記cDNA配列との間で所定の割合以上の類似性を有する部分を検索する検索手段と、前記ゲノム配列とcDNA配列をそれぞれグラフの縦軸と

横軸又は横軸と縦軸にとって前記検索手段で検索された部分を線分で該グラフ上に表示して前記 c D N A 配列に対応する前記ゲノム配列上の遺伝子のエクソン・イントロン構造を表示する表示手段とを備えることを特徴とするシーケンサ装置。

【請求項 6】 イントロン配列を跨ぐ相異なるエクソン領域内にあるプライマのペアを設計しこれを用いてゲノムライブラリと c D N A ライブラリとでそれぞれ P C R を行うステップと、該 P C R を行うステップで増幅されたゲノム配列と c D N A 配列とを入力するステップと、前記ゲノム配列中の所定の塩基長以上を有する部分配列について前記 c D N A 配列との間で所定の割合以上の類似性を有する部分を検索するステップと、前記ゲノム配列と c D N A 配列をそれぞれグラフの縦軸と横軸又は横軸と縦軸にとって前記検索するステップで検索された部分を線分で該グラフ上に表示することによってイントロン配列が存在するために異なったポリヌクレオチドが増幅されたことを表示して増幅したゲノム配列がイントロン配列を含んでいることを確認するステップと、を備えることを特徴とするプライマ設計方法。

【発明の詳細な説明】

【 0 0 0 1 】

【発明の属する技術分野】

本発明は遺伝子配列の情報解析に係わり、c D N A とゲノムの配列類似性検索結果から、ゲノム上の遺伝子の位置と構造を推定し表示する方法に関する。

【 0 0 0 2 】

【従来の技術】

ゲノム上の遺伝子の位置とそのエクソン・イントロン構造を推定する方法としては、c D N A 配列とゲノム配列との類似性検索を行い、類似性のある部分配列区間を列挙する方法がある。このとき、類似性のある部分配列区間は、類似度が高い順にソートして列挙される。類似度は、そのような類似性が偶然現れる確率によって評価し、その確率値が小さいものほど類似度が高いとする。

【 0 0 0 3 】

このようなソート法が有用である理由は、以下のように考えられる。生物のゲ

ノムは、遺伝子のコピーを派生させ分化させることにより進化してきた。そのため、一般に、一つの cDNA 配列に対して、ゲノム上の複数箇所に、種々の類似度で類似する部分配列が存在する。それら複数のゲノム部分配列のうち、実際にその cDNA の鋳型となった mRNA に転写されたゲノム部分配列は、類似度が最も高いものに限られる。このときの不一致部分は、SNP などの多型に起因するか、または、シーケンシング・エラーによるものと考えられる。従って、類似性のある区間を類似性の高い順にソートして列挙することにより、その cDNA の鋳型となった mRNA に転写されたゲノム上の部分配列が上位に列挙され、cDNA 配列とゲノム配列との対応付けが容易になる。

## 【0004】

また、cDNA 配列とゲノム配列との対応においては、cDNA 配列全体が一本の配列としてゲノム内の部分配列に対応することは少なく、一般には、cDNA 配列は何本かの部分配列に分かれ、その各々がゲノム内の部分配列に対応する。このような対応が見られる理由は、ヒトを含む真核生物において、ゲノムから mRNA が合成される際、スプライシングと呼ばれる現象が起きることによる。cDNA とゲノム上で対応している各々の部分配列はエクソンとよばれる。cDNA 上ではエクソンは切れ目なく繋がっているが、ゲノム上ではイントロンとよばれる部分配列を挟んで繋がっている。cDNA 上のエクソンとゲノム上のエクソンの位置関係は次のいずれかになっている。

## 【0005】

- (1) cDNA 上の各エクソン配列とゲノム上の各エクソン配列はほぼ一致し（以下、これらは向きが同じという）、それらは同じ順番で並んでいる。
- (2) cDNA 上の各エクソン配列とゲノム上の各エクソン配列は互いにほぼ相補鎖の関係にあり（以下、これらは向きが逆という）、それらは互いに反対の順番で並んでいる。

## 【0006】

このようなエクソン・イントロン構造をもつ cDNA 配列とゲノム配列との対応の様子は、類似性のある区間の列挙だけでは把握できず、それら類似性のある区間の相互の位置を調べる必要がある。そのためには、ゲノム配列上の塩基位置

と cDNA 配列上の塩基位置を両軸にとった 2 次元プロットが役立つ。最も単純なプロット法として、ゲノム配列の x 塩基目と cDNA 配列の y 塩基目が同一の塩基であるとき、2 次元上の座標 (x, y) に点をプロットする方法 (ドットマトリクス法) がある (105 頁、Sequence Analysis Primer, M. Gribskov and J. Devereux, Oxford University Press, 1992 年)。この方法では、局所的に精緻な比較が可能となる。また、より大局的な対応関係を捉える方法として、ゲノム配列内と cDNA 配列内に一定塩基長のウィンドウをとり、これらのウィンドウ内の塩基配列が一定割合以上類似しているとき、ゲノム配列内のウィンドウ位置を x 軸に cDNA 配列内のウィンドウ位置を y 軸にとり、それらのウィンドウに対応する線分を 2 次元平面上にプロットする方法がある (108 頁、Sequence Analysis Primer, M. Gribskov and J. Devereux, Oxford University Press, 1992 年)。この方法では、一塩基ずつの比較ではなく、数塩基～数十塩基ずつの平均的な比較が行われるため、より長い配列同士の比較が可能になり、また、偶然生じ意味をもたない短い一致部分を排除できる。

## 【0007】

## 【発明が解決しようとする課題】

エクソン・イントロン構造をもつ cDNA 配列とゲノム配列との対応関係を、判り易くグラフィック表示する。ゲノム上には多数の遺伝子が存在する領域があり、多数の cDNA が対応付けられる (貼り付けられるとも言う) ことがあり、それらの位置関係はグラフィック表示することにより、視覚的に理解しやすくなる。

## 【0008】

また、遺伝子のエクソン・イントロン構造において、イントロン配列はエクソン配列に比較して極めて長いことがある。cDNA 配列の長さは概ね数百から数万塩基長程度であるが、ゲノム上の対応する遺伝子領域は百万塩基長のオーダーまで広がることもある。このように cDNA とゲノムとで対応させるべき配列の長さが 3 桁も異なる場合には、同じサイズのウィンドウを cDNA 配列内とゲノム配列内で移動して調べる従来の方法は非効率的となる。

## 【0009】



また、ゲノム上の広い範囲にわたって cDNA との類似配列の位置を表示する場合、真の対応関係に関与しない多数の類似配列が現われ、真の対応関係を 2 次元表示の中から拾い出すことを妨げる。そのようなものとして、短い類似配列や、類似度の低い類似配列、向きや順番が不整合の類似配列などが考えられる。そこで、これらの不要な類似配列を除去することが必要になる。

【0010】

【課題を解決するための手段】

本発明では、与えられた cDNA 配列とゲノム断片配列に対して、以下の処理ステップから構成される方法によって、それらの間のエクソン・イントロン構造をもった対応関係を表示する。

(1) 与えられた cDNA 配列を纏めて検索用にデータベース化しておき、与えられた各ゲノム断片配列ごとに、それを検索配列として cDNA 配列データベースに対して類似性検索を繰り返し行うステップ。

【0011】

(2) 互いに類似性がある cDNA とゲノムの部分配列のペアを列挙し、そのペアの特徴量として、①部分配列の塩基長、②類似度、③各部分配列がゲノム配列上または cDNA 配列上で並ぶ向きと順番、④ cDNA 部分配列が他のペアの cDNA 部分配列と共同して cDNA 配列全体を被覆できる割合、を計算するステップ。

【0012】

(3) 前項で列挙された類似性のある部分配列ペアの集合の中から、上記の特徴量に関する所定の緩い条件を満たさないものを削除するステップ。これは、意味のある類似性を反映している可能性が低いものを除去して処理量を圧縮することを目的とする。即ち、所定の長さや所定の類似度に満たないもの、また、ゲノム上で互いに整合性のある向きと順番をとりえないもの、また、共同して cDNA 配列の所定以上の割合を覆う可能性のないものを除去する。

【0013】

(4) 前ステップで選び出された類似性のある部分配列ペアの集合の中から、上記の特徴量に関して更に厳格な条件により、表示すべきペアの集合を絞り込むス

テップ。これは、意味のある類似性を反映している可能性が高いものを正確に選  
び出すことを目的とする。そのためには、例えば、グラフィック表示を利用し、  
ユーザからの対話的な指示により絞込みの条件の閾値を与えるパラメータを調整  
する。または、ゲノム上に互いに整合性のある向きと順番で現われ、共同して c  
DNA 配列の所定以上の割合を覆うことができる部分配列の集合を、プログラム  
に従って自動的に選び出し、結果をグラフィック表示する。

## 【 0 0 1 4 】

(5) 選び出された cDNA とゲノムの部分配列ペアの位置関係を 2 次元的に表  
示するステップ。グラフの 1 の軸にゲノム配列上の塩基位置を、他の軸に cDN  
A 配列上の塩基位置をとり、各部分配列ペアを一本の線分で表示する。この線分  
は、それぞれ軸へ射影したときに部分配列の位置を表し、かつ、cDNA とゲノ  
ムの向きの対応を表す。

## 【 0 0 1 5 】

このため、本発明の cDNA 配列とゲノム配列との対応表示方法は、グラフの  
1 の軸にゲノム配列上の塩基位置を、他の軸に cDNA 配列上の塩基位置をとり  
、前記ゲノム配列のうち所定塩基長以上を有する部分配列について、前記 cDN  
A 配列との間で所定の割合以上の類似性を有する部分をグラフ上に線分で表示す  
ることを特徴とする。

また、複数の cDNA を縦軸にとり、前記 cDNA との対応関係を cDNA ご  
とに異なる色で表示することが好ましい。

## 【 0 0 1 6 】

また、本発明は、ゲノム配列と cDNA 配列とを入力するステップと、前記ゲ  
ノム配列中の所定の塩基長以上を有する部分配列について前記 cDNA 配列との  
間で所定の割合以上の類似性を有する部分を検索するステップと、前記ゲノム配  
列と cDNA 配列をそれぞれグラフの縦軸と横軸又は横軸と縦軸にとって前記検  
索するステップで検索された部分を線分で該グラフ上に表示するステップと、を  
備える cDNA 配列とゲノム配列との対応表示方法をコンピュータに実行させる  
ためのプログラムを記録したコンピュータ読み取り可能な記録媒体である。

## 【 0 0 1 7 】

さらに cDNA 配列とゲノム配列との対応表示方法は、前記所定の塩基長及び前記類似性の所定の割合を入力するステップを備えることが好ましい。

また、本発明のシーケンサ装置は、ネットワーク接続された又は内蔵されたゲノムデータベースにアクセスし、ゲノム配列を入力し、シーケンシングによって得られた cDNA 配列を入力する入力手段と、前記ゲノム配列中の所定の塩基長以上を有する部分配列について前記 cDNA 配列との間で所定の割合以上の類似性を有する部分を検索する検索手段と、前記ゲノム配列と cDNA 配列をそれぞれグラフの縦軸と横軸又は横軸と縦軸にとって前記検索手段で検索された部分を線分で該グラフ上に表示して前記 cDNA 配列に対応する前記ゲノム配列上の遺伝子のエクソン・イントロン構造を表示する表示手段とを備えることを特徴とする。

#### 【0018】

また、本発明のプライマ設計方法は、イントロン配列を跨ぐ相異なるエクソン領域内にあるプライマのペアを設計しこれを用いてゲノムライブラリと cDNA ライブラリとでそれぞれ PCR を行うステップと、該 PCR を行うステップで増幅されたゲノム配列と cDNA 配列とを入力するステップと、前記ゲノム配列中の所定の塩基長以上を有する部分配列について前記 cDNA 配列との間で所定の割合以上の類似性を有する部分を検索するステップと、前記ゲノム配列と cDNA 配列をそれぞれグラフの縦軸と横軸又は横軸と縦軸にとって前記検索するステップで検索された部分を線分で該グラフ上に表示することによってイントロン配列が存在するために異なったポリヌクレオチドが増幅されたことを表示して増幅したゲノム配列がイントロン配列を含んでいることを確認するステップと、を備えることを特徴とする。

#### 【0019】

##### 【発明の実施の形態】

以下、本発明の実施の形態を、図を用いて詳細に説明する。

図 1 に、与えられた cDNA 配列をデータベース内のゲノム配列に貼り付けることにより、cDNA に対応する遺伝子のエクソン・イントロン構造を可視化することを目的とした、本発明の一実施例における処理の流れを示す。

## 【0020】

図1において、101は解析の対象とするcDNA配列データであり、102はcDNA配列と比較されるべきゲノム配列を格納したデータベースである。103は、cDNA配列データとゲノム配列データベースを読み込む入力処理である。104は、以後の類似性検索に備えるために、入力されたcDNA配列データをデータベース化する処理であり、公知の方法を用いたプログラムformatdb (Altschul, Stephen F., Thomas L. Madden, Alejandro A. Schaffer, Jinghui Zhang, Zheng Zhang, Webb Miller, and David J. Lipman (1997), "Gapped BLAST and PSI-BLAST: a new generation of protein database search programs", Nucleic Acids Res. 25:3389-3402.) を使う。105は、ゲノムデータベース内にある各々のゲノム断片配列ごとに、それを検索配列としてcDNAデータベースに対して類似性検索処理を繰り返す処理である。この各々の類似性検索処理は、公知のアルゴリズムを用いたプログラムであるBLAST (Altschul, Stephen F., Thomas L. Madden, Alejandro A. Schaffer, Jinghui Zhang, Zheng Zhang, Webb Miller, and David J. Lipman (1997), "Gapped BLAST and PSI-BLAST: a new generation of protein database search programs", Nucleic Acids Res. 25:3389-3402.) を用いる。106は、各ゲノム断片配列ごとに得られた類似性検索結果を記述したテキストデータを全て読み込み、その中に現われる類似性がある部分配列を抽出して列挙し、その各々の部分配列を特徴付ける諸量を計算する処理である。107は、それらの諸特徴量に基づき、列挙された類似性のある部分配列の中から、所定の緩い条件を満たすものを選び出す、類似部分配列の1次選択処理である。これは、意味のある類似性を反映している可能性が低いものを除去して処理量を圧縮することを目的とする。その選択結果をファイル108に記憶する。ここまでの計算処理は時間を要するため、また、ここまでの計算は以後のユーザとの対話的処理とは独立に1回だけ行えばよいと、このようにファイルに記憶しておく。109は、cDNA上とゲノム上にある類似性のある部分配列のうちで選択されたものの相互の位置関係をファイル108から読み込んで、ユーザに分かり易く呈示するために、2次元のグラフィック表示データを生成する処理である。110はモニターディスプレイ、キーボード、マウスを備えたユーザインターフェース装置

で、109で生成されたグラフィックデータを表示するとともに、ユーザからの描画パラメータを受け付けて109に渡してグラフィックデータを再計算させ、109と110は共同して対話的な表示を行う。さらに、111は、更に厳格な条件により部分配列を更に絞り込む、類似部分配列の2次選択処理である。これは、意味のある類似性を反映している可能性が高いものをより正確に選び出すことを目的とする。110は、そのために必要となるパラメータをユーザから受け付け、それらを111に送る。111により更に絞り込まれた類似部分配列のデータは109に送られ、そこで、グラフィックデータが再計算される。これは、再び、110に送られ、ユーザに呈示される。109と110と111とにより、対話的に部分配列の選択の仕方を対話的に変更することができ、これにより、ゲノムとcDNAとの対応関係を正しく捉えた部分配列の集合を選び出すことができる。

#### 【0021】

図2は、106において、類似性のあるゲノム断片配列の部分配列とcDNA部分配列とのペアを抽出して得られるデータのデータ構造を表す。ここに現われる情報は全て、105のBLAST プログラムによる類似性検索結果の中から得ることができる。201は、1本のゲノム断片配列に対応するデータであり、全体のデータはこの繰り返し構造をもつ。201は、少なくとも、ゲノム断片配列を識別する名前とその配列長、及び、そのゲノム断片配列と類似性のある部分配列を持つcDNAに関する情報202の繰り返し構造を含む。202は、少なくとも、cDNAを識別する名前とその配列長、及び、ゲノムとの類似性のある部分配列に関する情報203の繰り返し構造を含む。以後、説明の簡略化のため、ゲノム内とcDNA内にある互いに類似性のある部分配列を“エクソン”と呼ぶことにする。これは生物学的なエクソンに対応することもあるが、それ以外に、偶然生じた類似性による部分配列のペアも含むこともある。203はエクソンの情報であり、少なくとも、長さ、ゲノムとcDNAとの一致塩基数、ゲノム断片配列内とcDNA配列内での位置の情報を含む。

#### 【0022】

図2に示したデータ構造は、図1内の106以降で処理される情報の基本構造であり、ファイル108に格納される情報もこのデータ構造をもつ。これは、106で得

られた情報から、107において有用性が低いと判断される一部の情報が除去されたものである。109は、図2に示したデータ構造をもつ情報を読み込んでグラフィック表示を行い、また、111は図2に示したデータ構造をもつ情報を読み込んで、そこから有用性が高いと判断されるエクソンを選び出し、再び、図2のデータ構造の情報を109に返す。

### 【0023】

図3は、107の類似部分配列ペア（エクソン）の1次選択処理の動作を説明するためのフローチャートである。301の終了判定を含む繰り返し処理により、全てのゲノム断片配列に対して以下の処理を行う。302で、処理中のゲノム断片配列に対する201に示す情報を読み込む。この中には、202に示すcDNAの情報が複数含まれる。303の終了判定を含む繰り返し処理により、これら全てのcDNAに対して以下の処理を行う。304で、処理中のcDNA配列に対する202に示す情報を読み込む。この中には、203に示すエクソンの情報が複数含まれる。305では、これらの個々のエクソンについて、

$$(\text{類似度}) = (\text{エクソン内一致塩基数}) / (\text{エクソン塩基長})$$

により類似度を計算し、これが所定の類似度に満たない場合は、203に列挙された中から該当エクソンを削除する。所定の類似度として、例えば80%を設定しておけば、現在処理中のcDNAの鋳型となった遺伝子（またはその近縁の遺伝子）に含まれるエクソン以外のゲノム断片の部分配列は、ほぼ除去されと考えられる。次に、306では、残ったエクソン長の最大値を求め、それが所定の値以上かどうかを判定する。多くの場合、遺伝子中のエクソンの中には100塩基長程度のもの少なくとも1つはある。したがって、例えば50塩基長程度の長さのエクソンがひとつも見つからないとすれば、この場合、ゲノム中に豊富に遍在する繰り返し配列の一部を捉えている可能性が高いと考えられるので、307によりすべてのエクソン情報とそのcDNA情報を除去する。そうでない場合は、エクソン長の合計を計算し、cDNA配列の全長との比を求め、308でその値が所定の値以上かを判定する。その比の値が例えば30%に満たないような場合は、それらのエクソンはcDNA配列のごく一部しか覆うことができないため、そこでのcDNAとゲノムとの関連は薄いと考えられるので、307によりすべてのエクソン情報

とその cDNA 情報を除去する。

#### 【0024】

図4は、109の表示処理により生成され、110のモニター画面上に描画されるイメージを、簡略化して表した説明図である。401は処理したゲノム断片配列のリストであり、その中の1項目（図では「ゲノム断片配列2」）が選択され、その項目に対する解析結果がモニター画面に表示されていることを表している。402は、横軸にゲノム断片配列上の塩基位置を荒い座標系（図ではメガ塩基単位）でとり、縦軸にcDNA配列上の塩基位置を細かい座標系（図ではキロ塩基単位）でとり、ゲノムとcDNA間の類似部分配列のペアを示すエクソンを線分で表す。これらのエクソンを表す線分は、実際のモニター画面では、cDNAごとに色分けして表示する。403は、各cDNAに対してエクソンの合併がcDNA配列の全体をどの割合まで覆うか示す。これは、そのcDNAが現在処理中のゲノム断片配列とどの程度強い関連があるかを表している。404はcDNA配列のリストであり、その中の1項目（図では「cDNA配列1」）が選択され、その項目に対する解析結果がモニター画面に表示されていることを表している。405は、404において選ばれたcDNAに対して、それを含む402のプロットの一部を拡大表示したものである。406は、405のエクソンを表す線分のプロットを、縦軸に射影したものである。ここで、エクソンの合併がcDNA全体をどの程度覆うかを確認できる。また、407は、405のエクソンを表す線分のプロットを、横軸に射影したものである。ここで、射影されたエクソンに挟まれた部分がイントロンを表す。408は、各エクソンに対して、その塩基長とその中の（ゲノム・cDNA間の）一致塩基数を表示したものである。これにより、各エクソンにおけるゲノム・cDNA間の類似度がどの程度高いかを確認できる。

#### 【0025】

図5は、111の類似部分配列ペア（エクソン）の2次選択処理の動作を説明するためのフローチャートである。501の終了判定を含む繰り返し処理により、全てのゲノム断片配列に対して以下の処理を行う。502で、処理中のゲノム断片配列に対する201に示す情報を読み込む。この中には、202に示すcDNAの情報が複数含まれる。503の終了判定を含む繰り返し処理により、これら全てのcDNA

Aに対して以下の処理を行う。504で、処理中のcDNA配列に対する202に示す情報を読み込む。この中には、203に示すエクソンの情報が複数含まれる。505では、これらの個々のエクソンについて、

$$(\text{類似度}) = (\text{エクソン内一致塩基数}) / (\text{エクソン塩基長})$$

により類似度を計算し、これが所望の類似度に満たない場合は、203に列挙された中から該当エクソンを削除する。所望の類似度は、ユーザインターフェース111によりプログラムに伝えられる。例えば、ここで類似度98%を要求すれば、2%程度の違いはSNPなどの多型またはシーケンシング・エラーによるものと許容して、現在処理中のcDNAの鋳型となった遺伝子（またはそれに酷似した遺伝子）に含まれるエクソンのみが選ばれると考えられる。次に、506では、残ったエクソンの集合を、向きと順番が互いに整合的であるようなグループに分割する。すなわち、各グループごとに、そこに属するエクソンの集合は次のいずれかの条件を満たす。

#### 【0026】

(1) cDNA上の各エクソン配列とゲノム上の各エクソン配列はほぼ一致し（これらは向きが同じ、または、正の向きという）、それらは同じ順番で並んでいる。

(2) cDNA上の各エクソン配列とゲノム上の各エクソン配列は互いにほぼ相補鎖の関係にあり（これらは向きが逆、または、負の向きという）、それらは互いに反対の順番で並んでいる。このようなグループ分けを行う手順は後述する。507の終了判定を含む繰り返し処理により、エクソンの各グループに対して以下の処理を行う。508でグループ内に属するエクソンの合併がcDNA全体を覆う割合を計算しそれが所定の割合（例えば95%）以上かを判定し、また、グループ内のエクソンをcDNA配列上で昇順に並べたとき隣り合うエクソン間の間隔が所定の塩基長（例えば10塩基）未満になっているかを判定し、違反があれば509においてこのグループに属する全エクソンを203から削除する。

#### 【0027】

1つのcDNAに属するエクソン全体を、506において上記のようにグループ分けするには、次のような手順に従う。まず、1つのcDNAに属するエクソン



全体を正・負の向きによって2つに分ける。次に、正の向きのエクソンをゲノム断片配列上の位置により昇順にソートし、また、負の向きのエクソンをゲノム断片配列上の位置により降順にソートする。それぞれの向きのエクソンについてソート順に見ていき、

(1) 最初のエクソンは新たなグループに属する。

【0028】

(2) 現在のエクソン q が直前に見たエクソン p に対して、

(q 右端塩基の cDNA 配列上での位置)

> (p 右端塩基の cDNA 配列上での位置) - (許容重なり塩基数)

が成り立つならば q は p と同じグループに属し、そうでない場合、q は新たなグループに属する。許容重なり塩基数としては、例えば、5 塩基程度でよい。

【0029】

【発明の実施の形態 - その2】

上記実施例による cDNA 配列とゲノム配列との対応表示を利用して、プライマ設計を行うための、本発明の第2の実施形態を、図を用いて詳細に説明する。

一般に、cDNA ライブラリを作成したとき、そこに含まれるポリヌクレオチドとして、cDNA 以外に、その他のゲノムの断片が紛れ込むことがある。従って、PCR を用いて cDNA 配列の一部を増幅しようとする際には、それが実際に cDNA 配列の一部であってそれ以外のゲノム断片ではないことを確認できることが有用である。

上記実施例を利用してプライマを設計することにより、このような確認が可能になる。

【0030】

図6は、そのようなプライマ設計法を説明する原理図である。601はゲノム上の塩基位置を表す軸であり、602はcDNA上の塩基位置を表す軸であり、603と604は一つのcDNAに属する相異なるエクソンを表す。603と604の塩基配列の中から、公知の方法(田平、林、PCR, PCR-SSCP法、新遺伝子工学ハンドブック、村松・山本編、75頁、羊土社、1999年)によりプライマ配列を選び出す。このプライマ配列のオリゴヌクレオチドを合成して、cDNAライブラリに対し

てPCRを行えば、これらのプライマは607、608の位置でcDNAに結合し、それらに挟まれた609に示すcDNAの部分配列をもつポリヌクレオチドが増幅される。一方、これと同じプライマを用いて、ゲノムライブラリに対してPCRを行えば、これらのプライマは610、611の位置でゲノムに結合し、それらに挟まれた612に示すゲノムの部分配列をもつポリヌクレオチドが増幅される。このポリヌクレオチドはイントロン配列を含んでいる。従って、これら2種類のPCRで増幅されたポリヌクレオチドの長さは異なる。

#### 【0031】

これに対して、cDNAライブラリの中に紛れ込んだゲノム断片からプライマを設計してしまった場合は、上記のような2種類のPCRで増幅されたポリヌクレオチドは一致する。651はゲノム上の塩基位置を表す軸であり、652はcDNA上の塩基位置を表す軸であり、653はエクソンを表す。653の塩基配列の中からプライマ配列を選び出す。このプライマ配列のオリゴヌクレオチドを合成して、cDNAライブラリに対してPCRを行えば、これらのプライマは656、657の位置でcDNAライブラリに含まれるゲノム断片に結合し、それらに挟まれた658に示す部分配列をもつポリヌクレオチドが増幅される。また、これと同じプライマを用いて、ゲノムライブラリに対してPCRを行えば、これらのプライマは659、660の位置でゲノムに結合し、それらに挟まれた661に示す配列をもつポリヌクレオチドが増幅される。これら2種類のPCRで増幅されたポリヌクレオチドは一致する。

#### 【0032】

このように、同じプライマを用いてcDNAライブラリとゲノムライブラリに対してPCRで増幅されたポリヌクレオチドの違いを調べることにより、cDNAに紛れ込んだゲノム断片ではなくcDNAの一部を増幅していることが確認できる。

#### 【0033】

##### 【発明の効果】

エクソン・イントロン構造をもつcDNA配列とゲノム配列との対応関係を、向きと順番が整合的な（エクソンに対応する）線分の並びとして、判り易くグラ

フィック表示する。エクソンの候補となる類似部分配列のペアについて、その両端の塩基位置と類似度等をあらかじめ計算しておき、その中からよりエクソンとして確からしい類似部分配列のペアを対話的に選んで描画するため、ゲノム上の広範囲にわたって高速に描画できる。短い類似配列や、類似度の低い類似配列、向きや順番が不整合の類似配列などを自動的に除去して表示するため、cDNA配列とゲノム配列との間の意味のある対応関係のみが描画される。

【図面の簡単な説明】

【図 1】

本発明の一実施の形態における処理の流れを示す図。

【図 2】

類似部分配列ペア（エクソン）を集めた情報のデータ構造。

【図 3】

類似部分配列ペア（エクソン）の 1 次選択処理の動作を説明するためのフローチャート。

【図 4】

モニター画面上に描画されるイメージを、簡略化して表した説明図。

【図 5】

類似部分配列ペア（エクソン）の 2 次選択処理の動作を説明するためのフローチャート。

【図 6】

本発明の第 2 の実施形態におけるプライマ設計法の原理を説明する図

【符号の説明】

- 101 解析の対象とする cDNA 配列データ
- 102 cDNA 配列と比較されるべきゲノム配列を格納したデータベース
- 103 cDNA 配列データとゲノム配列データベースを読み込む入力処理
- 104 類似性検索処理のために cDNA 配列データをデータベース化する処理
- 105 各ゲノム断片配列を検索配列として cDNA データベースに対して類似性検索を繰り返す処理
- 106 類似性がある部分配列ペア（エクソン）を抽出してその特徴量を計算す

る処理

107 処理量を圧縮することを目的とした、類似部分配列の 1 次選択処理

108 類似性がある部分配列ペア（エクソン）によりゲノムと cDNA を対応  
付けるデータを格納したファイル

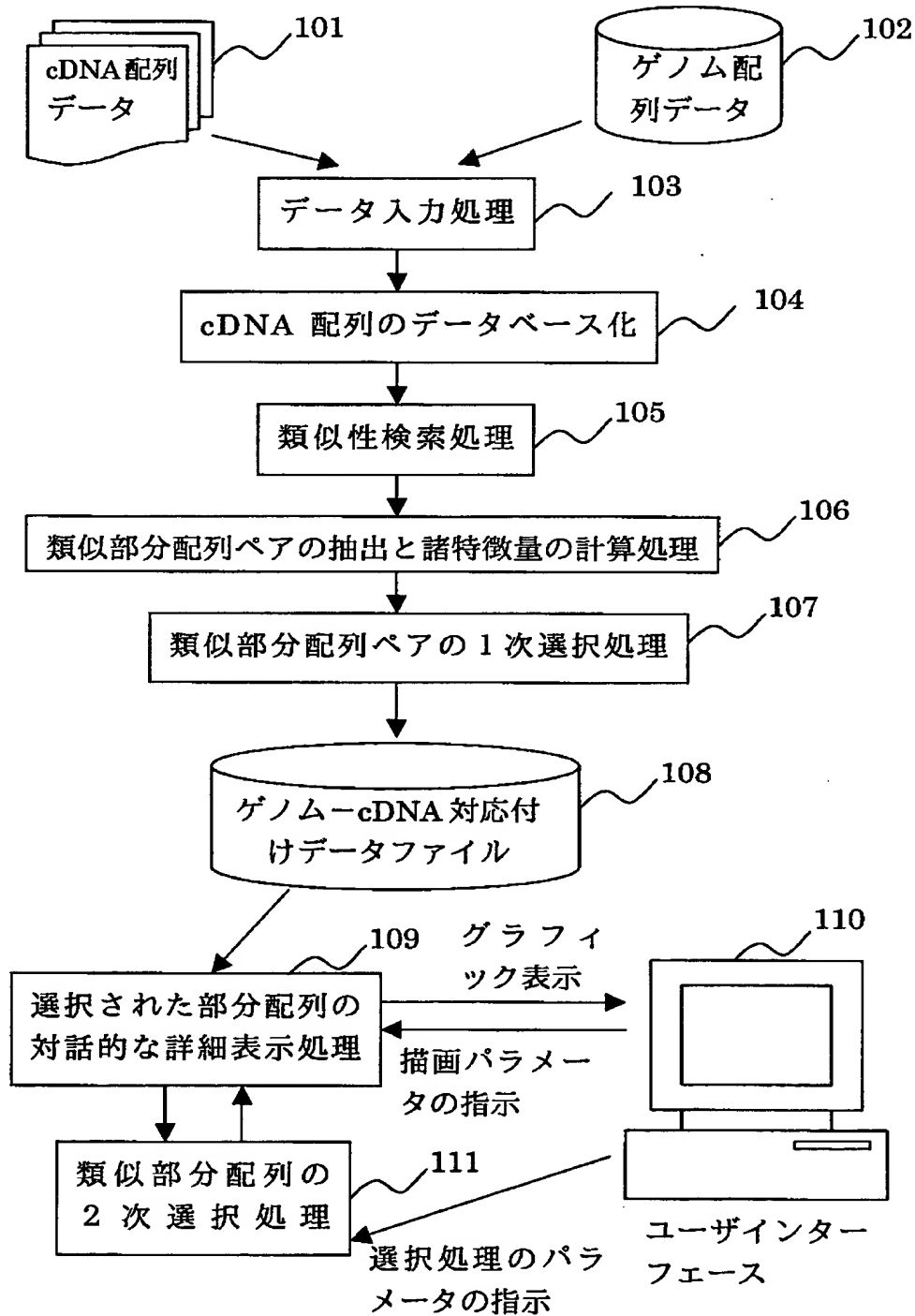
109 2 次元のグラフィック表示データを生成する処理

110 ユーザインターフェース装置

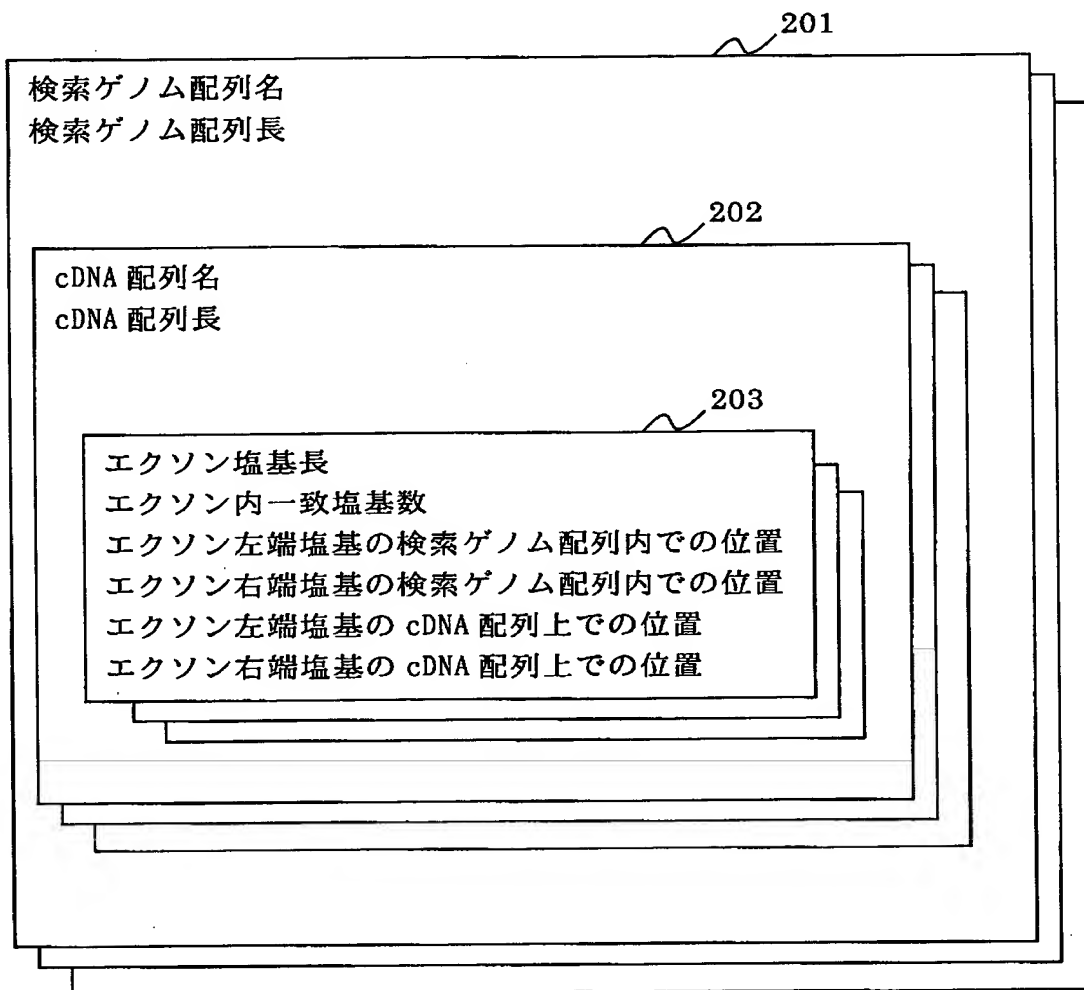
111 意味のある類似部分配列を正確に選び出すための 2 次選択処理

【書類名】 図面

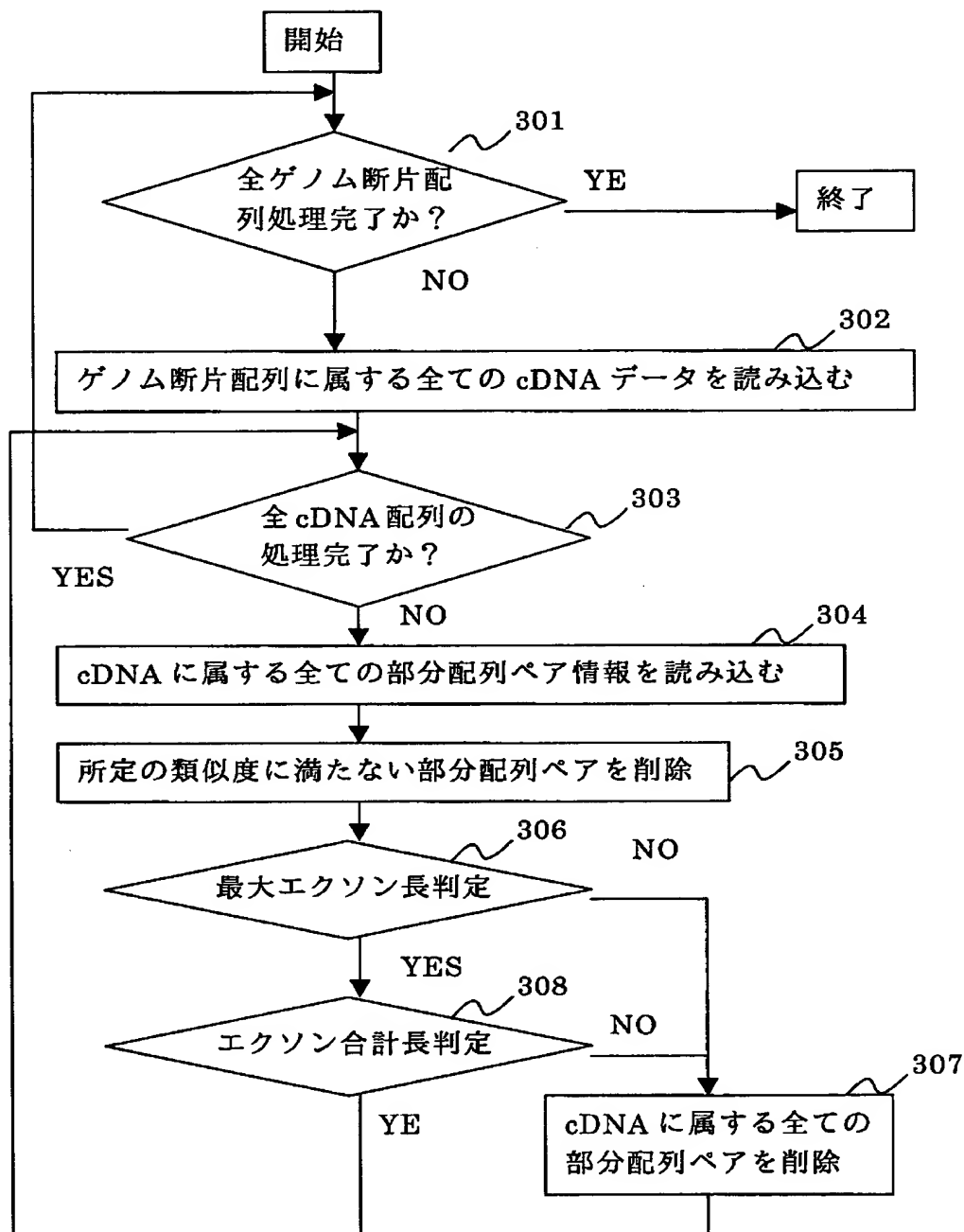
【図 1】



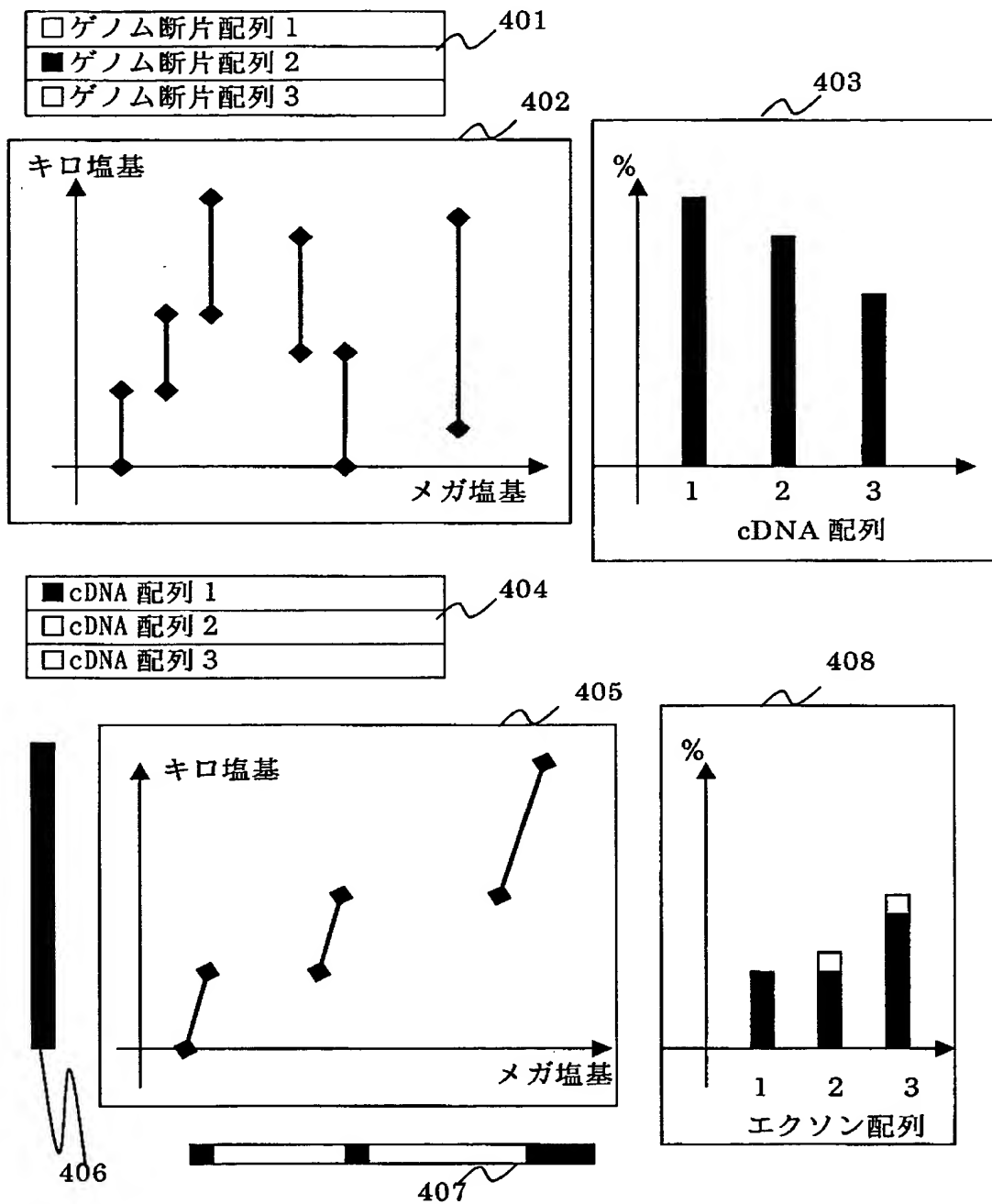
【図 2】



【図 3】

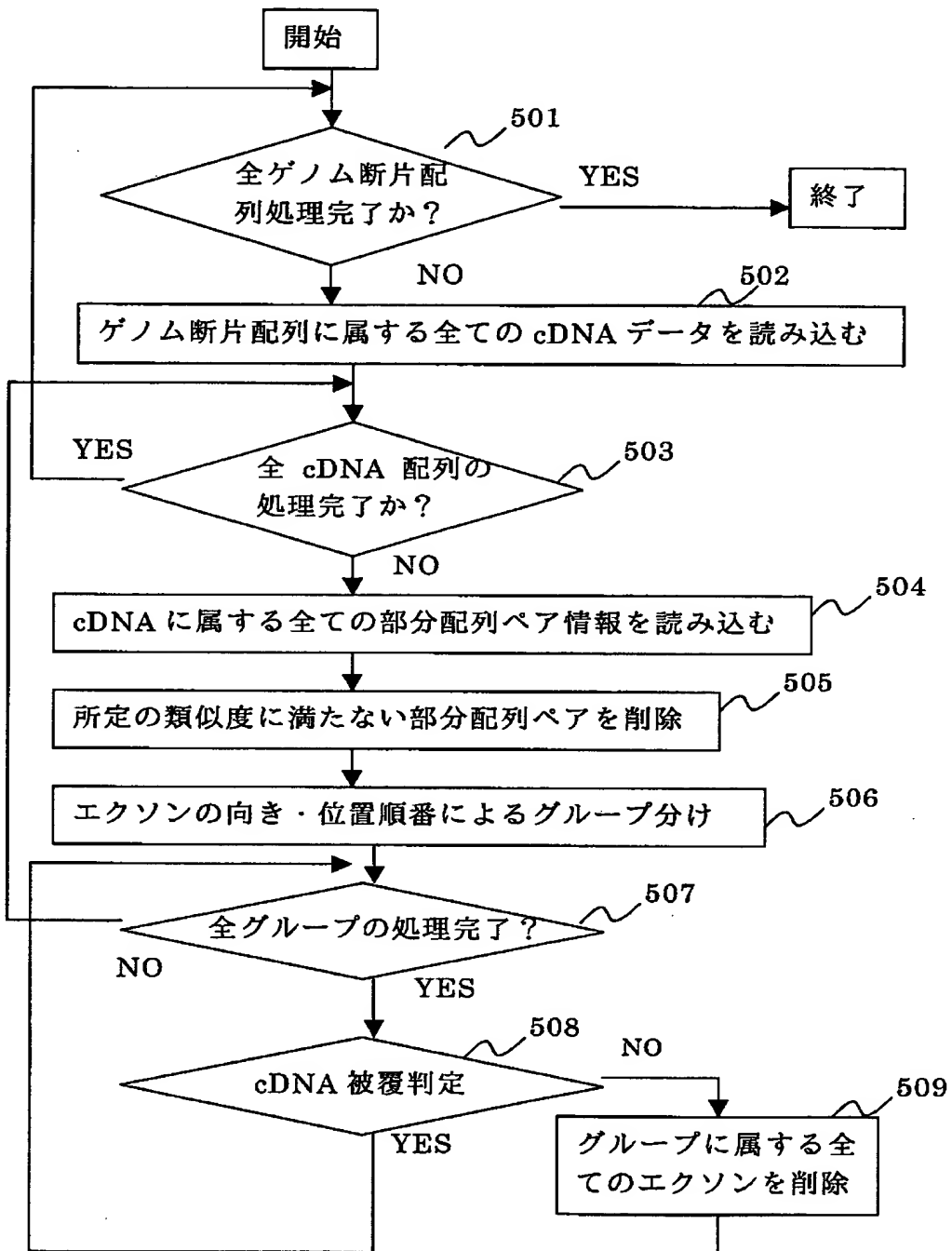


【図 4】

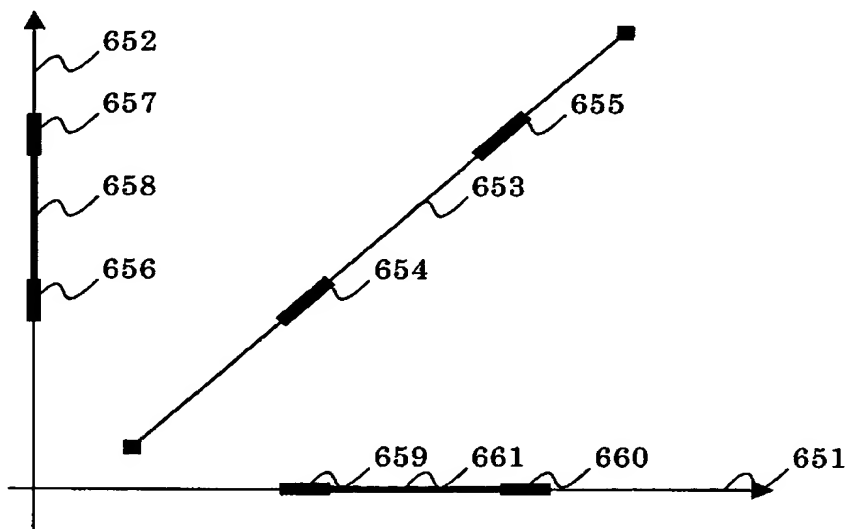
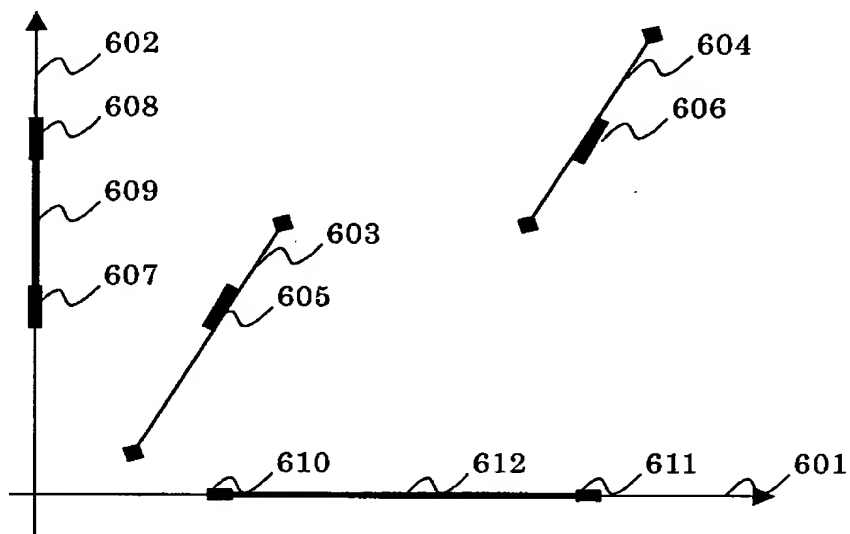




【図 5】



【図 6】



【書類名】 要約書

【要約】

【課題】 エクソン・イントロン構造をもつ cDNA 配列とゲノム配列との対応関係を、判り易くグラフィック表示する表示方法を提供する。

【解決手段】 cDNA とゲノムとの類似性検索結果から、類似性のある部分配列のペア（エクソン）の両端の塩基位置・類似度等の情報を抽出する。これらの中から、類似度・配列長などの観点で意味をもつ可能性が低いと判定される部分配列ペアの情報を除去する。また、エクソン間の向き・順番に関する整合性を調べ、cDNA を所定の割合以上被覆して cDNA との関連が明確になるエクソンのみを選択する。選択されたエクソンを、グラフの 1 の軸にゲノム配列上の塩基位置を、他の軸に cDNA 配列上の塩基位置をとり、線分で表示することにより、線分の並びとしてイントロン・エクソン構造が視覚的に確認できる。

【選択図】 図 1

出 願 人 履 歴 情 報

識別番号 [000005108]

1. 変更年月日 1990年 8月31日

[変更理由] 新規登録

住 所 東京都千代田区神田駿河台4丁目6番地  
氏 名 株式会社日立製作所